## Solution 1: Lasso Regularization

(a) First of all, we will use the fact that **X** has orthonormal columns to show that:

$$\hat{\boldsymbol{\theta}} = \left(\underbrace{\mathbf{X}^T \mathbf{X}}_{I}\right)^{-1} \mathbf{X}^T \mathbf{y}$$

$$\hat{\boldsymbol{\theta}} = \mathbf{X}^T \mathbf{y}$$
(1)

We will now use the result of eq. 1 to show that:

$$\arg\min_{\boldsymbol{\theta}} \mathcal{R}_{reg}(\boldsymbol{\theta}) = \arg\min_{\boldsymbol{\theta}} \frac{1}{2} \sum_{i=1}^{n} \left( y^{(i)} - \boldsymbol{\theta}^{\top} \mathbf{X}^{(i)} \right)^{2} + \lambda \sum_{i=1}^{p} |\theta_{i}|$$

$$= \arg\min_{\boldsymbol{\theta}} \frac{1}{2} ||\mathbf{y} - \mathbf{X}\boldsymbol{\theta}||_{2}^{2} + \lambda ||\boldsymbol{\theta}||_{1}$$

$$= \arg\min_{\boldsymbol{\theta}} \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^{T} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \lambda ||\boldsymbol{\theta}||_{1}$$

$$= \arg\min_{\boldsymbol{\theta}} \frac{1}{2} \left( \underbrace{\mathbf{y}^{T} \mathbf{y}}_{indep.\ of\ \boldsymbol{\theta}} - 2 \underbrace{\mathbf{y}^{T} \mathbf{X}}_{\boldsymbol{\theta}^{T}} \boldsymbol{\theta} + \boldsymbol{\theta}^{T} \underbrace{\mathbf{X}^{T} \mathbf{X}}_{\boldsymbol{I}} \boldsymbol{\theta} \right) + \lambda ||\boldsymbol{\theta}||_{1}$$

$$= \arg\min_{\boldsymbol{\theta}} -\hat{\boldsymbol{\theta}}^{T} \boldsymbol{\theta} + \frac{\boldsymbol{\theta}^{T} \boldsymbol{\theta}}{2} + \lambda ||\boldsymbol{\theta}||_{1}$$

$$= \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^{p} -\hat{\theta}_{i} \theta_{i} + \frac{\theta_{i}^{2}}{2} + \lambda ||\boldsymbol{\theta}_{i}||.$$

$$(2)$$

- (b) The advantage of this representation if we are interested in finding  $\boldsymbol{\theta}$  is that we can optimize each  $g_i(\theta_i)$  separately to get optimal entries for  $\theta_1, \ldots, \theta_p$ .
- (c) Let's use the hint and compare  $g_i(\theta_i)$  with  $g_i(-\theta_i)$ , for the case where  $\theta_i > 0$ :

$$g_i(\theta_i) = -\hat{\theta}_i \theta_i + \frac{\theta_i^2}{2} + \lambda |\theta_i|$$

$$g_i(-\theta_i) = +\hat{\theta}_i \theta_i + \frac{(-\theta_i)^2}{2} + \lambda |\theta_i|$$
(3)

We also know that non-positive  $\theta_i$  have always a greater or equal value for  $g_i(\theta_i)$  than than their positive counterpart:

$$\theta_i > 0 \longrightarrow g_i(\theta_i) \le g_i(-\theta_i)$$
 (4)

The second and third term of both equations in eq. 3 are equivalent, so we can ignore them. Accordingly, to comply with the condition from eq. 4, the minimizer  $\theta_i^*$  must be non-negative.

$$-\hat{\theta}_i \theta_i < \hat{\theta}_i \theta_i \longrightarrow \theta_i^* > 0 \tag{5}$$

(d) To calculate the minimizer  $\theta_i^*$ , we will derive with respect to the parameter and set the derivative to zero:

$$\frac{\partial g_{i}(\theta_{i})}{\partial \theta_{i}} = -\hat{\theta}_{i} + \theta_{i} + \underbrace{\lambda}_{\theta_{i} > 0}$$

$$= -\hat{\theta}_{i} + \theta_{i} + \lambda \stackrel{!}{=} 0 \longrightarrow \theta_{i} = \underbrace{\hat{\theta}_{i} - \lambda}_{\geq 0}$$

$$\theta_{i}^{*} = \underbrace{\hat{\theta}_{i} - \lambda}_{\geq 0}$$

$$= \max(0, \hat{\theta}_{i} - \lambda)$$
(6)

(e) Let's use the hint again and compare  $g_i(\theta_i)$  with  $g_i(-\theta_i)$ , for the case where  $\theta_i < 0$ :

$$g_{i}(\theta_{i}) = -\hat{\theta}_{i}\theta_{i} + \frac{\theta_{i}^{2}}{2} + \lambda|\theta_{i}|$$

$$g_{i}(-\theta_{i}) = +\hat{\theta}_{i}\theta_{i} + \frac{(-\theta_{i})^{2}}{2} + \lambda|\theta_{i}|$$
(7)

$$\theta_i < 0 \longrightarrow g_i(\theta_i) \ge g_i(-\theta_i)$$
 (8)

The second and third term in both equations of 7 are equivalent. Using the condition from eq. 8, we can conclude that the minimizer  $\theta_i^*$  must be non-positive.

$$-\hat{\theta}_i \theta_i \ge \hat{\theta}_i \theta_i \longrightarrow \theta_i^* \le 0 \tag{9}$$

(f) Calculating the derivative again and setting it to zero, we get:

$$\frac{\partial g_i(\theta_i)}{\partial \theta_i} = -\hat{\theta}_i + \theta_i \underbrace{-\lambda}_{\theta_i < 0} \stackrel{!}{=} 0 \longrightarrow \theta_i = \underbrace{\hat{\theta}_i + \lambda}_{\leq 0}$$

$$\theta_i^* = \underbrace{\hat{\theta}_i + \lambda}_{\leq 0}$$

$$= \min(0, \hat{\theta}_i + \lambda)$$
(10)

(g)

$$\begin{cases}
\theta_{i} > 0 \longrightarrow \theta_{i}^{*} = \overbrace{1}^{\operatorname{sign}(\hat{\theta}_{i})} \cdot \max(0, \widehat{\theta}_{i}^{*} - \lambda) \\
\theta_{i} \leq 0 \longrightarrow \theta_{i}^{*} = \min(0, \widehat{\theta}_{i} + \lambda) = \underbrace{-1}_{\operatorname{sign}(\hat{\theta}_{i})} \cdot \max(0, -\widehat{\theta}_{i}^{*} - \lambda)
\end{cases} \tag{11}$$

By transforming the min function from 10 to a max function, we can combine the two cases in only one expression:

$$\theta_i^* = \operatorname{sign}(\hat{\theta}_i) \max(|\hat{\theta}_i| - \lambda, 0) \tag{12}$$